

برای انجام رگرسیون خطی از تابع $lm()$ استفاده می‌کنیم:

$lm(\text{formula}, \text{data}, \dots)$

formula: مدل رگرسیونی است که می‌خواهیم به داده‌ها برازش کنیم. به عنوان مثال به صورت

$Y \sim X_1 + X_2$ می‌نویسیم که در آن X_1 و X_2 متغیرهای ترفیعی در مقیاس استاندارد است

data: ظرف مجموعه داده‌هایی است که متغیرها در آن وارد دارند.

مثال: مجموعه داده‌های Highway1 موجود در R را در نظر بگیرید:

> $rm(\text{dist} = \text{ls}())$

> Highway1

> $\text{head}(\text{Highway1})$ → چندتا مشاهده از مجموعه داده را می‌بینیم

rate: نرخ تصادفات ۱۹۷۳ در ایالت‌های جنوبی آمریکا در هر مایل جاده

len: طول پرتگاه به مایل

adtb: متوسط تعداد اتومبیل روزانه به ۱۰۰۰

trks: نسبت کامیون‌ها به تعداد کل خودروها

sigsl: تعداد سیگنال‌ها در هر مایل جاده

slim: محدودیت سرعت

shld: زیبایی سطح جاده

lane: تعداد لاین‌های جاده

acpb: تعداد نقاط دسترسی در هر مایل

wt: تعداد تقسیمات نوع آزاد راه در هر مایل

lwid: عرض لاین‌ها

htype: نوع جاده

> $\text{my model} = \text{lm}(\text{rate} \sim \text{acpb} + \text{slim} + \text{len} + \text{shld}, \text{Highway1})$

> mymodel # or print(mymodel)

دل‌های بازج داده شده می‌باشند مقادیر تریسیدی، انجی‌دی و ...

$$\hat{y}_i = 9.17574 + 0.10136 \text{ acpt}_i + 0.09808 \text{ slm}_i - 0.07517 \text{ len}_i$$

$$-0.01099 \text{ shld}_i$$

> summary(mymodel)

خلاصه‌ای از نتایج بازج مدل به شکل آماره‌های آزمون‌های مدل. پارامترهای تریسیدی به همراه پارامترهای سایر متغیرها. آماره‌های آزمون t و F - مقادیر آزمون‌های مقادیر تریسیدی. همچنین R^2 ، R^2_{adj} و \sqrt{MSE} و جدول ANOVA.

$$\underline{y} = X \underline{\beta} + \underline{\epsilon}$$

$$y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_k x_{ik} + \epsilon_i \quad i=1, \dots, n$$

$$\underline{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} \quad X = \begin{bmatrix} 1 & x_{11} & \dots & x_{1k} \\ \vdots & \vdots & \dots & \vdots \\ 1 & x_{n1} & \dots & x_{nk} \end{bmatrix} \quad \underline{\beta} = \begin{bmatrix} \beta_0 \\ \vdots \\ \beta_k \end{bmatrix}$$

$$F = \frac{SSR/k}{MSE}$$

$$F > F_{1-\alpha}(k, n-p) \Rightarrow H_0$$

$$\begin{cases} H_0: \beta_1 = \dots = \beta_k = 0 \\ H_1: 0 \neq \dots \end{cases}$$

آزمون استاندارد برای مدل

$$\hat{\underline{\beta}} = (X'X)^{-1} X' \underline{y}$$

$$V(\hat{\underline{\beta}}) = \sigma^2 (X'X)^{-1} = \sigma^2 C$$

$$V(\hat{\beta}_j) = \sigma^2 C_{jj}$$

$$\widehat{V(\hat{\beta}_j)} = MSE C_{jj}$$

$$Cov(\hat{\beta}_j, \hat{\beta}_r) = \sigma^2 C_{jr}$$

$$Cov(\hat{\beta}_j, \hat{\beta}_r) = MSE C_{jr}$$

$$R^2 = \frac{SSR}{SSY} = 1 - \frac{SSE}{SSY}$$

$$MSE = \frac{SSE}{n-p}$$

$$R^2_{adj} = 1 - \frac{MSE}{MSY}$$

$$= 1 - \frac{SSE/n-p}{SSY/n-1}$$

$$\begin{matrix} \hat{\beta}_0 \\ \hat{\beta}_1 \\ \vdots \end{matrix} \quad \begin{matrix} \sqrt{MSE C_{00}} \\ \sqrt{MSE C_{11}} \\ \vdots \end{matrix} \quad \begin{matrix} T = \frac{\hat{\beta}_0 - 0}{\sqrt{MSE C_{00}}} \\ T = \frac{\hat{\beta}_1 - 0}{\sqrt{MSE C_{11}}} \\ \vdots \end{matrix}$$

$$p\text{-value} = P(|T| > |T_{\alpha}|)$$

$$p\text{-value} = P(|T| > |T_{\alpha}|)$$

$$\begin{cases} H_0: \beta_1 = 0 \\ H_1: \beta_1 \neq 0 \end{cases}$$

* آزمون H_0 را قبول کنیم به قدری که X_j در حضور ϵ مقادیر تریسیدی یا بازی نیست

$$T^2 = F = \frac{\hat{\beta}_1^2}{MSE C_{11}}$$

$$p\text{-value} = P(F > F_{\alpha})$$

> anova(myModel)

معیار خوبی قبل است با این معیار که آن‌ها می‌تواند از T یا F یا R^2 معیار R^2 زودتر R^2 را انتخاب می‌کند.

برای استناد $SSR = SSE + SSY$ به این ترتیب عمل می‌کنیم:

> SSE = anova(myModel)["Residuals", "sum Sq"]

> n = length(Highway1 & rate)

> SSY = t(Highway1 & rate) %*% Highway1 & rate - n *
(mean(Highway1 & rate))^2

> SSR = SSY - SSE

تست t ها z ها $\hat{\beta}_j$ ها به ترتیب داریم:

> myModel\$ fitted.values # or fitted.values(myModel) or
fitted(myModel)

> myModel\$residuals # or residuals(myModel) or
resid(myModel)

> myModel\$ coefficients # or coefficients(myModel)

برای استناد آن‌ها در فاصله‌های فیت و ریسونی:

> confint(myModel, level = 0.9)

$$\hat{\beta}_j \pm t_{(n-p)} \sqrt{MS E C_{jj}}$$

برای استناد آن‌ها در حالت $ln(\beta)$ داریم:

> vcov(myModel)

ماتریس واریانس کوواریانس $\hat{\beta}$ ها.

رسم نمودار خط روند با این داده‌ها:

در صورتی که مدل ln یک تقصیر تقویمی یا در تقصیر تقویمی باشد همان‌طور که در خط روند دیده می‌شود.

داده‌ها را رسم کرد.

> myModel1 = lm(rate ~ acpt, Highway1)

> plot(Highway1 & acpt, Highway & rate)

> abline(myModel1, col = "red", lwd = 2.)

نکته: ابتدا باید car و mpg را نصب کنید پس

> library(car)

> scatter3d(Highway1\$rate ~ Highway1\$acpt + Highway1\$slim,

ارتوسیم تقویرهای باند. $\left. \begin{array}{l} \text{revolutions} = 3, \text{ speed} = 0.5 \\ \text{grid} = F \end{array} \right\}$

چنانچه بخواهیم از تمامی مقیورهای موجود در مجموعه داده، به عنوان مقیور کوفتهی مدل استفاده کنیم:

> mymodel2 = lm(rate ~ ., Highway1)

> mymodel2
برای مقیورهای رتبه‌ای و مقیورهای نامرکز شده استفاده می‌شود.

آنچه بخواهیم رگرسیون عمودی از مرکز به ازای رسم:

> mymodel3 = lm(rate ~ acpt + slim + len + shld - 1,

Highway1)

می‌توانیم فواصل اطمینان برای میانگین‌های پاسخ و فواصل پیش‌بینی برای پاسخ‌ها رسم کنیم.

> predict(mymodel, newdata = list(acpt = 4.5, slim = 50, len = 6
, shld = 11), interval = "confidence", level = 0.95)

فواصل اطمینان میانگین پاسخ

$$\hat{y}_0 \pm t_{(n-p)} \sqrt{MSE \left(\frac{1}{n} + \underline{\underline{x_0' (X'X)^{-1} x_0}} \right)} \quad x_0 = \begin{bmatrix} 1 \\ x_{01} \\ \vdots \\ x_{0k} \end{bmatrix}$$

> predict(

, interval = "prediction", level = 0.95)
فواصل پیش‌بینی

$$\hat{y}_0 \pm t_{(n-p)} \sqrt{MSE \left(1 + \frac{1}{n} + x_0' (X'X)^{-1} x_0 \right)}$$

